

WHAT IS CLAIMED IS:

- 1 1. A method comprising:
 - 2 a) accepting a packet associated with a flow;
 - 3 b) generating a flow group identifier from the flow;
 - 4 c) determining whether any other packets associated with the
 - 5 flow group are present in a switch fabric;
 - 6 d) if it is determined that other packets associated with
 - 7 the flow group are present in the switch fabric, then
 - 8 assigning the packet to a path being used by the flow group,
 - 9 and if it is determined that other packets associated with
 - 10 the flow group are not present in the switch fabric, then
 - 11 assigning the packet to a path using path congestion status
 - 12 information.
- 1 2. The method of claim 1 wherein the act of generating a flow
- 2 group identifier from the flow includes hashing a flow identifier.
- 1 3. The method of claim 1 wherein the act of determining whether
- 2 any other packets associated with the flow group are present in a
- 3 switch fabric includes maintaining an outstanding packet counter.
- 1 4. The method of claim 3 wherein the outstanding packet counter
- 2 is associated with the flow group identifier.
- 1 5. The method of claim 4 wherein the act of maintaining an
- 2 outstanding packet counter includes incrementing the outstanding
- 3 packet counter each time a packet belonging to the flow group is
- 4 sent into the switch fabric, and decrementing the outstanding

5 packet counter each time a packet belonging to the flow group
6 leaves the switch fabric.

1 6. The method of claim 5 wherein the act of decrementing the
2 outstanding packet counter is performed in response to a message
3 from an output port.

1 7. The method of claim 6 further comprising:
2 - passing the message from the output port to a
3 corresponding input element,
4 - passing the message from the corresponding input element,
5 through the switch fabric, to another output element, and
6 - passing the message from the other output element to
7 another input element corresponding to the other output
8 element, wherein the other input element originated the
9 packet.

1 8. The method of claim 3 wherein the act of maintaining the
2 outstanding packet counter includes resetting the outstanding
3 packet counter if it remains non-zero for more than a
4 predetermined period of time.

1 9. The method of claim 1 wherein the act of assigning the packet
2 to a path using path congestion status information includes
3 - selecting a switch plane having at least one uncongested
4 path, and
5 - selecting an uncongested path of the selected switch
6 plane.

1 10. The method of claim 9 wherein the act of selecting a switch
2 plane having at least one uncongested path uses a round robin
3 discipline.

1 11. The method of claim 9 wherein the act of selecting an
2 uncongested path of the selected switch plane uses a round robin
3 discipline.

1 12. A machine-readable medium having stored thereon a data
2 structure comprising a plurality of entries, each entry including
3 a) a flow group identifier,
4 b) an outstanding packet in switch fabric indicator, and
5 c) a path identifier.

1 13. The machine-readable medium of claim 12 further including a
2 second data structure comprising a plurality of entries, each
3 entry including
4 a) the path identifier, and
5 b) path status information.

1 14. The machine-readable medium of claim 13 wherein the path
2 status information includes
3 i) an indicator of whether or not the path has failed,
4 and
5 ii) an indicator of whether or not the path is
6 congested.

1 15. Apparatus comprising:
2 a) an input for accepting a packet associated with a flow;
3 b) means for generating a flow group identifier from the
4 flow;
5 c) means for determining whether any other packets
6 associated with the flow group are present in a switch
7 fabric;

8 d) means for assigning the packet to a path being used by
9 the flow group if it is determined that other packets
10 associated with the flow group are present in the switch
11 fabric, and for assigning the packet to a path using path
12 congestion status information if it is determined that other
13 packets associated with the flow group are not present in the
14 switch fabric.

1 16. The apparatus of claim 15 wherein the means for generating a
2 flow group identifier from the flow hash a flow identifier.

1 17. The apparatus of claim 15 wherein the means for determining
2 whether any other packets associated with the flow group are
3 present in a switch fabric maintain an outstanding packet counter.

1 18. The apparatus of claim 17 wherein the outstanding packet
2 counter is associated with the flow group identifier.

1 19. The apparatus of claim 18 wherein the means for maintaining
2 an outstanding packet counter increment the outstanding packet
3 counter each time a packet belonging to the flow group is sent
4 into the switch fabric, and decrement the outstanding packet
5 counter each time a packet belonging to the flow group leaves the
6 switch fabric.

1 20. The apparatus of claim 19 wherein the decrementing of the
2 outstanding packet counter is performed in response to a message
3 from an output port.

1 21. The apparatus of claim 20 further comprising:
2 - means for passing the message from the output port to a
3 corresponding input element,

4 - means for passing the message from the corresponding input
5 element, through the switch fabric, to another output
6 element, and
7 - means for passing the message from the other output
8 element to another input element corresponding to the other
9 output element, wherein the other input element originated
10 the packet.

1 22. The apparatus of claim 17 wherein the means for maintaining
2 the outstanding packet counter reset the outstanding packet
3 counter if it remains non-zero for more than a predetermined
4 period of time.

1 23. The apparatus of claim 15 wherein the means for assigning the
2 packet to a path using path congestion status information include
3 means for
4 - selecting a switch plane having at least one uncongested
5 path, and
6 - selecting an uncongested path of the selected switch
7 plane.

1 24. The apparatus of claim 23 wherein the means for selecting a
2 switch plane having at least one uncongested path use a round
3 robin discipline.

1 25. The apparatus of claim 24 wherein the means for selecting an
2 uncongested path of the selected switch plane use a round robin
3 discipline.

1 26. A method for alleviating head-of-line blocking in an
2 input-buffered switch, wherein the switch includes a plurality of

input modules, each input module including virtual output queues and virtual path queues, the method comprising:

- a) assigning an incoming cell to an appropriate one of the virtual output queues using cell destination information;
- b) providing a head-of-line cell of the one of the virtual output queues to an appropriate one of the virtual path queues using path identifier information of the cell;
- c) for an input module-to-switch plane link, selecting one of a number of virtual path queues associated with the link and having at least one cell; and
- d) sending the cell from the selected one of the number of virtual path queues over the link.

27. The method of claim 26 wherein the path identifier information of the cell was provided using a dynamic hashing scheme.

28. The method of claim 26 further comprising:

- e) determining whether or not the cell sent over the link was the last cell of a packet; and
- f) if it was determined that the cell sent over the link was the last cell of a packet, then instructing the virtual output queue to send cells of a next packet to an appropriate one of the virtual path queues.

29. For use in a switch, an input module comprising:

- a) a plurality of virtual output queues for accepting cells; and
- b) a plurality of virtual path queues for accepting head-of-line cells from the plurality of virtual output queues.

1 30. The input module of claim 29 wherein the number of the
2 plurality of virtual output queues equals a number of output ports
3 of the switch.

1 31. The input module of claim 29 wherein the number of the
2 plurality of virtual path queues equals a number of paths through
3 a switch fabric of the switch.

1 32. The input module of claim 29 wherein the number of the
2 plurality of virtual path queues equals a product of (a) a number
3 of switch planes of a switch fabric of the switch and (b) a number
4 of paths through each of the switch planes.